

# acmqueue Thinking Clearly about Performance

**Improving the performance of complex software is difficult, but understanding some fundamental principles can make it easier.**

**Cary Millsap, Method R Corporation**

When I joined Oracle Corporation in 1989, performance—what everyone called “Oracle tuning”—was difficult. Only a few people claimed they could do it very well, and those people commanded high consulting rates. When circumstances thrust me into the “Oracle tuning” arena, I was quite unprepared. Recently, I’ve been introduced to the world of “MySQL tuning,” and the situation seems very similar to what I saw in Oracle more than 20 years ago.

It reminds me a lot of how difficult beginning algebra seemed when I was about 13 years old. At that age, I had to appeal heavily to trial and error to get through. I can remember looking at an equation such as  $3x + 4 = 13$  and basically stumbling upon the answer,  $x = 3$ .

The trial-and-error method worked—albeit slowly and uncomfortably—for easy equations, but it didn’t scale as the problems got tougher—for example,  $3x + 4 = 14$ . Now what? My problem was that I wasn’t thinking clearly yet about algebra. My introduction at age 15 to teacher James R. Harkey put me on the road to solving that problem.

In high school Mr. Harkey taught us what he called an *axiomatic* approach to solving algebraic equations. He showed us a set of steps that worked every time (and he gave us plenty of homework to practice on). In addition, by executing those steps, we necessarily *documented* our thinking as we worked. Not only were we thinking clearly, using a reliable and repeatable sequence of steps, but we were also *proving* to anyone who read our work that we were thinking clearly.

Our work for Mr. Harkey looked like this:

$3.1x + 4$	$= 13$	problem statement
$3.1x + 4 - 4$	$= 13 - 4$	subtraction property of equality
$3.1x$	$= 9$	additive inverse property, simplification
$3.1x / 3.1$	$= 9 / 3.1$	division property of equality
$x$	$\approx 2.903$	multiplicative inverse property, simplification

This was Mr. Harkey’s axiomatic approach to algebra, geometry, trigonometry, and calculus: one small, logical, provable, and auditable step at a time. It’s the first time I ever really *got* mathematics.

Naturally, I didn’t realize it at the time, but of course *proving* was a skill that would be vital for my success in the world after school. In life I’ve found that, of course, knowing things matters, but *proving* those things to other people matters more. Without good *proving* skills, it’s difficult to be a good consultant, a good leader, or even a good employee.

My goal since the mid-1990s has been to create a similarly rigorous approach to Oracle performance optimization. Lately, I have been expanding the scope of that goal beyond Oracle to: “Create an axiomatic approach to computer software performance optimization.” I’ve found that not

many people like it when I talk like that, so let's say it like this: "My goal is to help you think clearly about how to optimize the performance of your computer software."

#### WHAT IS PERFORMANCE?

Googling the word *performance* results in more than a half-billion hits on concepts ranging from bicycle racing to the dreaded employee review process that many companies these days are learning to avoid. Most of the top hits relate to the subject of this article: the time it takes for computer software to perform whatever task you ask it to do.

And that's a great place to begin: the *task*, a business-oriented unit of work. Tasks can nest: "print invoices" is a task; "print one invoice"—a subtask—is also a task. For a computer user, *performance* usually means the time it takes for the system to execute some task. *Response time* is the execution duration of a task, measured in time per task, such as "seconds per click." For example, my Google search for the word *performance* had a response time of 0.24 seconds. The Google Web page rendered that measurement right in my browser. That is evidence to me that Google values my perception of Google performance.

Some people are interested in another performance measure: *throughput*, the count of task executions that complete within a specified time interval, such as "clicks per second." In general, people who are responsible for the performance of *groups* of people worry more about throughput than does the person who works in a solo contributor role. For example, an individual accountant is usually more concerned about whether the response time of a daily report will require that accountant to stay late after work. The manager of a group of accounts is also concerned about whether the system is capable of processing all the data that all of the accountants in that group will be processing.

#### RESPONSE TIME VERSUS THROUGHPUT

Throughput and response time have a generally reciprocal type of relationship, but not exactly. The real relationship is subtly complex.

**EXAMPLE 1.** Imagine that you have measured your throughput at 1,000 tasks per second for some benchmark. What, then, is your users' average response time? It's tempting to say that the average response time is  $1/1,000 = .001$  seconds per task, but it's not necessarily so.

Imagine that the system processing this throughput had 1,000 parallel, independent, homogeneous service channels (that is, it's a system with 1,000 independent, equally competent service providers, each awaiting your request for service). In this case, it is possible that each request consumed exactly 1 second.

Now, you can know that average response time was somewhere between 0 and 1 second per task. You cannot derive response time exclusively from a throughput measurement, however; you have to measure it separately (I carefully include the word *exclusively* in this statement, because there are mathematical models that can compute response time for a given throughput, but the models require more input than just throughput).

The subtlety works in the other direction, too. You can certainly flip this example around and prove it. A scarier example, however, will be more fun.

**EXAMPLE 2.** Your client requires a new task that you're programming to deliver a throughput of 100 tasks per second on a single-CPU computer. Imagine that the new task you've written executes

in just .001 seconds on the client's system. Will your program yield the throughput that the client requires?

It's tempting to say that if you can run the task once in just a thousandth of a second, then surely you'll be able to run that task at least 100 times in the span of a full second. And you're right, if the task requests are nicely serialized, for example, so that your program can process all 100 of the client's required task executions inside a loop, one after the other.

But what if the 100 tasks per second come at your system at random, from 100 different users logged into your client's single-CPU computer? Then the gruesome realities of CPU schedulers and serialized resources (such as Oracle latches and locks and writable access to buffers in memory) may restrict your throughput to quantities much less than the required 100 tasks per second. It might work; it might not. You cannot derive throughput exclusively from a response time measurement. You have to measure it separately.

Response time and throughput are not necessarily reciprocals. To know them both, you need to measure them both. Which is more important? For a given situation, you might answer legitimately in either direction. In many circumstances, the answer is that both are vital measurements requiring management. For example, a system owner may have a business requirement not only that response time must be 1.0 second or less for a given task in 99 percent or more of executions but also that the system must support a sustained throughput of 1,000 executions of the task within a 10-minute interval.

#### PERCENTILE SPECIFICATIONS

In the prior section, I used the phrase "in 99 percent or more of executions" to qualify a response time expectation. Many people are more accustomed to such statements as "average response time must be  $r$  seconds or less." The percentile way of stating requirements maps better, though, to the human experience.

**TABLE 1**

**Two Lists of Response Times  
(average for each is 1.000 second)**

	List A	List B
1	.924	.796
2	.928	.798
3	.954	.802
4	.957	.823
5	.961	.919
6	.965	.977
7	.972	1.076
8	.979	1.216
9	.987	1.273
10	1.373	1.320

**EXAMPLE 3.** Imagine that your response time tolerance is 1 second for some task that you execute on your computer every day. Imagine further that the lists of numbers shown in table 1 represent the measured response times of 10 executions of that task. The average response time for each list is 1.000 second. Which one do you think you would like better?

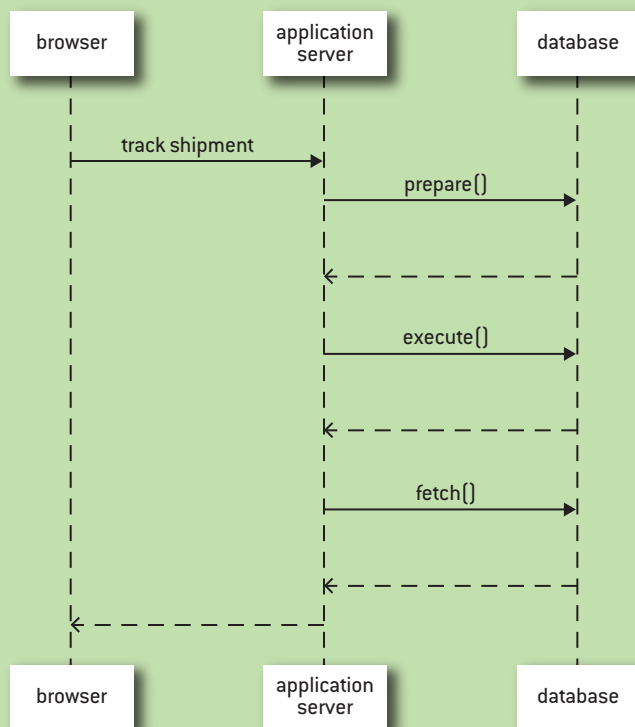
Although the two lists in table 1 have the same average response time, the lists are quite different in character. In list A, 90 percent of response times were one second or less. In list B, only 60 percent of response times were one second or less. Stated in the opposite way, list B represents a set of user experiences of which 40 percent were dissatisfactory, but list A (having the same average response time as list B) represents only a 10 percent dissatisfaction rate.

In list A, the 90th percentile response time is .987 seconds; in list B, it is 1.273 seconds. These statements about percentiles are more informative than merely saying that each list represents an average response time of 1.000 second.

As GE says, “Our customers feel the variance, not the mean.”<sup>4</sup> Expressing response-time goals as percentiles makes for much more compelling requirement specifications that match with end-user expectations: for example, the “Track Shipment” task must complete in less than .5 second in at least 99.9 percent of executions.

## FIGURE 1

**UML Sequence Diagram Showing Interactions among a Browser, an Application Server, and a Database**



## PROBLEM DIAGNOSIS

In nearly every performance problem I've been invited to repair, the stated problem has been about response time: "It used to take less than a second to do X; now it sometimes takes 20+." Of course, a specific statement like that is often buried under veneers of other problems such as: "Our whole [adjectives deleted] system is so slow we can't use it."<sup>8</sup>

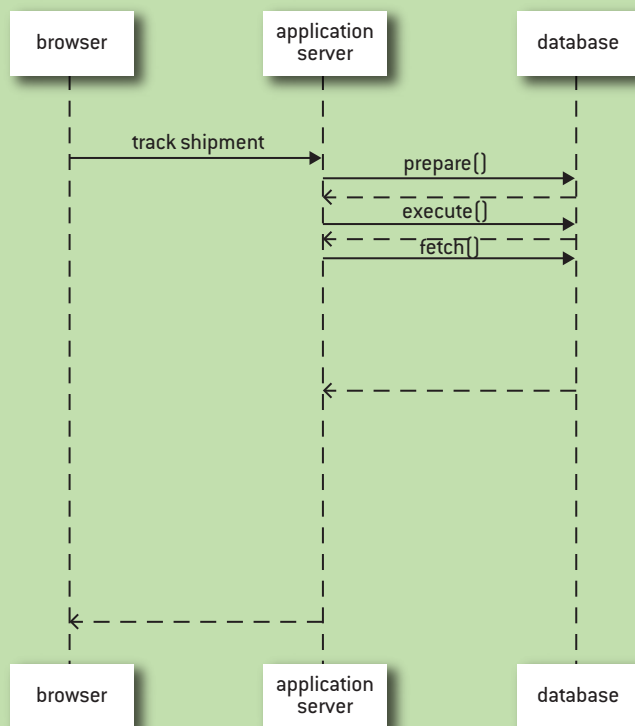
Just because something has happened a lot for me doesn't mean that it will happen for you. The most important thing for you to do first is state the problem clearly, so that you can then think about it clearly.

A good way to begin is to ask, what is the goal state that you want to achieve? Find some specifics that you can *measure* to express this: for example, "Response time of X is more than 20 seconds in many cases. We'll be happy when response time is one second or less in at least 95 percent of executions." That sounds good in theory, but what if your user doesn't have such a specific quantitative goal? This particular goal has two quantities (1 and 95); what if your user doesn't know either one of them? Worse yet, what if your user *does* have specific ideas, but those expectations are impossible to meet? How would you know what "possible" or "impossible" even is?

Let's work our way up to those questions.

## FIGURE 2

**UML Sequence Diagram Drawn to Scale, Showing the Response Time Consumed at Each Tier in the System**



## THE SEQUENCE DIAGRAM

A sequence diagram is a type of graph specified in UML (Unified Modeling Language), used to show the interactions between objects in the sequential order that those interactions occur. The sequence diagram is an exceptionally useful tool for visualizing response time. Figure 1 shows a standard UML sequence diagram for a simple application system composed of a browser, application server, and a database.

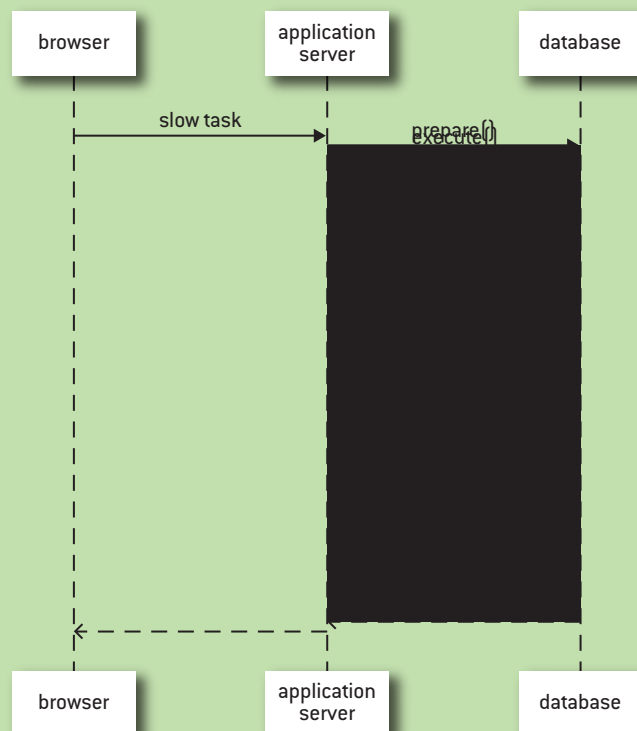
Imagine now drawing the sequence diagram to scale, so that the distance between each “request” arrow coming in and its corresponding “response” arrow going out is proportional to the duration spent servicing the request. I’ve shown such a diagram in Figure 2. This is a good graphical representation of how the components represented in your diagram are spending your user’s time. You can “feel” the relative contribution to response time by looking at the picture.

Sequence diagrams are just right for helping people conceptualize how their responses are consumed on a given system, as one tier hands control of the task to the next. Sequence diagrams also work well to show how simultaneous processing threads work in parallel, and they are good tools for analyzing performance outside of the information technology business.<sup>10</sup>

The sequence diagram is a good conceptual tool for talking about performance, but to think

## FIGURE 3

**UML Sequence Diagram Showing 322,968 Database Calls**



clearly about performance, you need something else. Here's the problem. Imagine that the task you're supposed to fix has a response time of 2,468 seconds (41 minutes, 8 seconds). In that period of time, running that task causes your application server to execute 322,968 database calls. Figure 3 shows what the sequence diagram for that task would look like.

There are so many request and response arrows between the application and database tiers that you can't see any of the detail. Printing the sequence diagram on a very long scroll isn't a useful solution, because it would take weeks of visual inspection before you would be able to derive useful information from the details you would see.

The sequence diagram is a good tool for conceptualizing flow of control and the corresponding flow of time. To think clearly about response time, however, you need something else.

## THE PROFILE

The sequence diagram doesn't scale well. To deal with tasks that have huge call counts, you need a convenient aggregation of the sequence diagram so that you understand the most important patterns in how your time has been spent. Table 2 shows an example of a *profile*, which does the trick. A *profile* is a tabular decomposition of response time, typically listed in descending order of component response time contribution.

**EXAMPLE 4.** The profile in table 2 is rudimentary, but it shows exactly where your slow task has spent your user's 2,468 seconds. With the data shown here, for example, you can derive the percentage of response time contribution for each of the function calls identified in the profile. You can also derive the average response time for each type of function call during your task.

A profile shows *where your code has spent your time* and—sometimes even more importantly—where it has *not*. There is tremendous value in not having to guess about these things.

From the data shown in table 2, you *know* that 70.8 percent of your user's response time is consumed by `DB:fetch()` calls. Furthermore, if you can drill down into the individual calls whose durations were aggregated to create this profile, you can know how many of those `App:await_db_netIO()` calls corresponded to `DB:fetch()` calls, and you can know how much response time each of those consumed.

## TABLE 2

**Profile Showing the Decomposition of a 2,468.000-second Response Time**

	Function call	R (sec)	Calls
1	DB: fetch()	1,748.229	322,968
2	App: await_db_netIO()	338.470	322,968
3	DB: execute()	152.654	39,142
4	DB: prepare()	97.855	39,142
5	Other	58.147	89,422
6	App: render_graph()	48.274	7
7	App: tabularize()	23.481	4
8	App: read()	0.890	2
	<b>Total</b>	<b>2,468.000</b>	

With a profile, you can begin to formulate the answer to the question, “How long *should* this task run?” which, by now, you know is an important question in the first step of any good problem diagnosis.

#### AMDAHL'S LAW

Profiling helps you think clearly about performance. Even if Gene Amdahl hadn't given us Amdahl's law back in 1967, you would probably have come up with it yourself after the first few profiles you looked at.

Amdahl's law states: Performance improvement is proportional to how much a program uses the thing you improved. If the thing you're trying to improve contributes only 5 percent to your task's total response time, then the maximum impact you'll be able to make is 5 percent of your total response time. This means that the closer to the top of a profile that you work (assuming that the profile is sorted in descending response-time order), the bigger the benefit potential for your overall response time.

This doesn't mean that you always work through a profile in top-down order, though, because you also need to consider the *cost* of the remedies you'll be executing.<sup>9</sup>

**EXAMPLE 5.** Consider the profile in Table 3. It's the same profile as in Table 2, except here you can see how much time you think you can save by implementing the best remedy for each row in the profile, and you can see how much you think each remedy will cost to implement.

Which remedy would you implement first? Amdahl's law says that implementing the repair on line 1 has the greatest potential benefit of saving about 851 seconds (34.5 percent of 2,468 seconds). If it is truly “super expensive,” however, then the remedy on line 2 may yield better net benefit—and that's the constraint to which you really need to optimize—even though the potential for response time savings is only about 305 seconds.

A tremendous value of the profile is that you can learn exactly how much improvement you should expect for a proposed investment. It opens the door to making much better decisions about what remedies to implement first. Your predictions give you a yardstick for measuring your own

### TABLE 3

**Profile Showing the Potential Improvement and Corresponding Cost (difficulty) for Each Table 2 Line Item**

	Potential improvement % and cost of investment	R (sec)	R (%)
1	34.5% super expensive	1,748.229	70.8%
2	12.3% dirt cheap	338.470	13.7%
3	Impossible to improve	152.654	6.2%
4	4.0% dirt cheap	97.855	4.0%
5	0.1% super expensive	58.147	2.4%
6	1.6% dirt cheap	48.274	2.0%
7	Impossible to improve	23.481	1.0%
8	0.0% dirt cheap	0.890	0.0%
	<b>Total</b>	<b>2,468.000</b>	



performance as an analyst. Finally, it gives you a chance to showcase your cleverness and intimacy with your technology as you find more efficient remedies for reducing response time more than expected, at lower-than-expected costs.

What remedy action you implement first really boils down to how much you trust your cost estimates. Does “dirt cheap” really take into account the risks that the proposed improvement may inflict upon the system? For example, it may seem dirt cheap to change that parameter or drop that index, but does that change potentially disrupt the good performance behavior of something out there that you’re not even thinking about right now? Reliable cost estimation is another area in which your technological skills pay off.

Another factor worth considering is the political capital that you can earn by creating small victories. Maybe cheap, low-risk improvements won’t amount to much overall response-time improvement, but there’s value in establishing a track record of small improvements that exactly fulfill your predictions about how much response time you’ll save for the slow task. A track record of prediction and fulfillment ultimately—especially in the area of software performance, where myth and superstition have reigned at many locations for decades—gives you the credibility you need to influence your colleagues (your peers, your managers, your customers...) to let you perform increasingly expensive remedies that may produce bigger payoffs for the business.

A word of caution, however: don’t get careless as you rack up successes and propose ever-bigger, costlier, riskier remedies. Credibility is fragile. It takes a lot of work to build it up but only one careless mistake to bring it down.

## SKEW

When you work with profiles, you repeatedly encounter sub-problems such as this:

**EXAMPLE 6.** The profile in table 2 revealed that 322,968 DB: fetch() calls had consumed 1,748.229 seconds of response time. How much unwanted response time would be eliminated if you could eliminate half of those calls? The answer is almost never, “Half of the response time.” Consider this far simpler example for a moment:

**EXAMPLE 7.** Four calls to a subroutine consumed four seconds. How much unwanted response time would be eliminated if you could eliminate half of those calls? The answer depends upon the response times of the individual calls that we could eliminate. You might have assumed that each of the call durations was the average  $4/4 = 1$  second, but nowhere did the statement tell you that the call durations were uniform.

**EXAMPLE 8.** Imagine the following two possibilities, where each list represents the response times of the four subroutine calls:

A = {1, 1, 1, 1}

B = {3.7, .1, .1, .1}

In list A, the response times are uniform, so no matter which half (two) of the calls you eliminate, you will reduce total response time to two seconds. In list B, however, it makes a tremendous difference which two calls are eliminated. If you eliminate the first two calls, then the total response time will drop to .2 seconds (a 95 percent reduction). If you eliminate the final two calls, then the total response time will drop to 3.8 seconds (only a 5 percent reduction).

Skew is a nonuniformity in a list of values. The possibility of skew is what prohibits you from providing a precise answer to the question I asked at the beginning of this section. Let's look again:

**EXAMPLE 9.** The profile in table 2 revealed that 322,968 DB: fetch() calls had consumed 1,748.229 seconds of response time. How much unwanted response time would you eliminate by eliminating half of those calls? Without knowing anything about skew, the most precise answer you can provide is, "Somewhere between 0 and 1,748.229 seconds."

Imagine, however, that you had the additional information available in table 4. Then you could formulate much more precise best-case and worst-case estimates. Specifically, if you had information like this, you would be smart to try to figure out how to specifically eliminate the 47,444 calls with response times in the .01- to .1-second range.

#### MINIMIZING RISK

A couple of sections back I mentioned the risk that repairing the performance of one task can damage the performance of another. This reminds me of something that happened to me once in Denmark. It's a quick story:

**SCENE:** The kitchen table in Måløv, Denmark; *the* oak table, in fact, of Oak Table Network fame, a network of Oracle practitioners who believe in using scientific methods to improve the development and administration of Oracle-based systems.<sup>12</sup> Roughly 10 people sit around the table, working on their laptops and conducting various conversations.

**CARY:** Guys, I'm burning up. Would you mind if I opened the window for a little bit to let some cold air in?

**CAREL-JAN:** Why don't you just take off your heavy sweater?

**THE END.**

There's a general principle at work here that humans who optimize know: when everyone is happy except for you, make sure your local stuff is in order before you go messing around with the global stuff that affects everyone else, too.

## TABLE 4

**Skew Histogram for the 322,968 Calls from Table 2**

	Range {min ≤ e < max}		R (sec)	Calls
1	0	.000001	.000	0
2	.000001	.00001	.002	397
3	.00001	.0001	.141	2,169
4	.0001	.001	31.654	92,557
5	.001	.01	746.852	180,399
6	.01	.1	1,688.451	47,444
7	.1	1	.900	2
<b>Total</b>			<b>2,468.000</b>	<b>322,968</b>

This principle is why I flinch whenever someone proposes to change a system's Oracle SQL\*Net packet size when the problem is really a couple of badly written Java programs that make unnecessarily many database calls (and, hence, unnecessarily many network I/O calls as well). If everybody is getting along fine except for the user of one or two programs, then the safest solution to the problem is a change whose scope is localized to just those one or two programs.

## EFFICIENCY

Even if everyone on the entire system is suffering, you should still focus first on the program that the business needs fixed first. The way to begin is to ensure that the program is working as efficiently as it can. *Efficiency* is the inverse of how much of a task execution's total service time can be eliminated without adding capacity and without sacrificing required business function.

In other words, efficiency is an inverse measure of waste. Here are some examples of waste that commonly occur in the database application world:

- A middle-tier program creates a distinct SQL statement for every row it inserts into the database. It executes 10,000 database prepare calls (and thus 10,000 network I/O calls) when it could have accomplished the job with one prepare call (and thus 9,999 fewer network I/O calls).
- A middle-tier program makes 100 database fetch calls (and thus 100 network I/O calls) to fetch 994 rows. It could have fetched 994 rows in 10 fetch calls (and thus 90 fewer network I/O calls).
- A SQL statement (my choice of wording here is a dead giveaway that I was introduced to SQL within the Oracle community) touches the database buffer cache 7,428,322 times to return a 698-row result set. An extra filter predicate could have returned the seven rows that the end user really wanted to see, with only 52 touches upon the database buffer cache.

Certainly, if a system has some global problem that creates inefficiency for broad groups of tasks across the system (e.g., ill-conceived index, badly set parameter, poorly configured hardware), then you should fix it. Don't tune a system to accommodate programs that are inefficient, however. (Admittedly, sometimes you need a tourniquet to keep from bleeding to death, but don't use a stopgap measure as a permanent solution. Address the inefficiency.) There is a lot more leverage in curing the program inefficiencies themselves. Even if the programs are commercial off-the-shelf applications, it will benefit you more in the long run to work with your software vendor to make your programs efficient than it will to try to optimize your system to be as efficient as it can with an inherently inefficient workload.

Improvements that make your program more efficient can produce tremendous benefits for everyone on the system. It's easy to see how top-line reduction of waste helps the response time of the task being repaired. What many people don't understand as well is that making one program more efficient creates a side effect of performance improvement for other programs on the system that have no apparent relation to the program being repaired. It happens because of the influence of *load* upon the system.

## LOAD

*Load* is competition for a resource induced by concurrent task executions. It is the reason that the performance testing done by software developers doesn't catch all the performance problems that show up later in production.

One measure of load is *utilization*, which is resource usage divided by resource capacity for a

specified time interval. As utilization for a resource goes up, so does the response time a user will experience when requesting service from that resource. Anyone who has ridden in an automobile in a big city during rush hour has experienced this phenomenon. When the traffic is heavily congested, you have to wait longer at the tollbooth.

The software you use doesn't actually "go slower" as your car does when you're going 30 mph in heavy traffic instead of 60 mph on the open road. Computer software always goes the same speed no matter what (a constant number of instructions per clock cycle), but certainly response time degrades as resources on your system get busier.

There are two reasons that systems get slower as load increases: *queuing delay* and *coherency delay*.

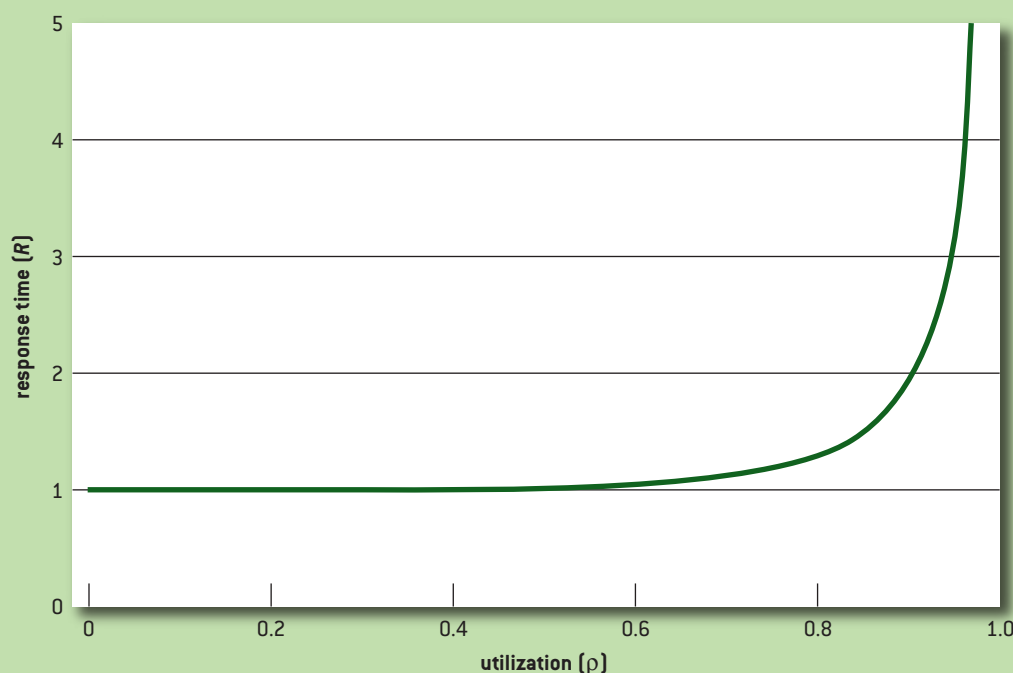
#### QUEUING DELAY

The mathematical relationship between load and response time is well known. One mathematical model, called  $M/M/m$ , relates response time to load in systems that meet one particularly useful set of specific requirements.<sup>11</sup> One of the assumptions of  $M/M/m$  is that the system you are modeling has "theoretically perfect scalability." This is akin to having a physical system with "no friction," an assumption that so many problems in introductory physics courses invoke.

Regardless of some overreaching assumptions such as the one about perfect scalability,  $M/M/m$  has a lot to teach us about performance. Figure 4 shows the relationship between response time and load using  $M/M/m$ .

## FIGURE 4

**Curve Showing Response Time as a Function of Utilization for an  $M/M/m$  System with  $m = 8$  Service Channels**



In figure 4 you can see mathematically what you feel when you use a system under different load conditions. At low load, your response time is essentially the same as your response time at no load. As load ramps up, you sense a slight, gradual degradation in response time. That gradual degradation doesn't really do much harm, but as load continues to ramp up, response time begins to degrade in a manner that's neither slight nor gradual. Rather, the degradation becomes quite unpleasant and, in fact, hyperbolic.

*Response time (R)*, in the perfect scalability M/M/m model, consists of two components: *service time (S)* and *queuing delay (Q)*, or  $R = S + Q$ . *Service time* is the duration that a task spends consuming a given resource, measured in time per task execution, as in *seconds per click*. *Queuing delay* is the time that a task spends enqueued at a given resource, awaiting its opportunity to consume that resource. Queuing delay is also measured in time per task execution (e.g., seconds per click).

In other words, when you order lunch at Taco Tico, your response time (*R*) for getting your order is the queuing delay time (*Q*) that you spend in front of the counter waiting for someone to take your order, plus the service time (*S*) you spend waiting for your order to hit your hands once you begin talking to the order clerk. Queuing delay is the difference between your response time for a given task and the response time for that same task on an otherwise unloaded system (don't forget our *perfect scalability* assumption).

## THE KNEE

When it comes to performance, you want two things from a system:

- The best response time you can get: you don't want to have to wait too long for tasks to get done.
- The best throughput you can get: you want to be able to cram as much load as you possibly can onto the system so that as many people as possible can run their tasks at the same time.

Unfortunately, these two goals are contradictory. Optimizing to the first goal requires you to minimize the load on your system; optimizing to the second goal requires you to maximize it. You can't do both simultaneously. Somewhere in between—at some load level (that is, at some utilization value)—is the optimal load for the system.

The utilization value at which this optimal balance occurs is called the *knee*. This is the point at which throughput is maximized with minimal negative impact to response times. (I am engaged in an ongoing debate about whether it is appropriate to use the term *knee* in this context. For the time being, I shall continue to use it. See sidebar for details.) Mathematically, the knee is the utilization value at which response time divided by utilization is at its minimum. One nice property of the knee is that it occurs at the utilization value where a line through the origin is tangent to the response-time curve. On a carefully produced M/M/m graph, you can locate the knee quite nicely with just a straightedge, as shown in figure 5.

Another nice property of the M/M/m knee is that you need to know the value of only one parameter to compute it. That parameter is the number of parallel, homogeneous, independent *service channels*. A *service channel* is a resource that shares a single queue with other identical resources, such as a booth in a toll plaza or a CPU in an SMP (symmetric multiprocessing) computer.

The italicized lowercase *m* in the term M/M/*m* is the number of service channels in the system being modeled. The M/M/*m* knee value for an arbitrary system is difficult to calculate, but I've done it for you in table 5, which shows the knee values for some common service channel counts. (By this point, you may be wondering what the other two Ms stand for in the M/M/*m* queuing model name.

They relate to assumptions about the randomness of the timing of your incoming requests and the randomness of your service times. See [http://en.wikipedia.org/wiki/Kendall%27s\\_notation](http://en.wikipedia.org/wiki/Kendall%27s_notation) for more information, or *Optimizing Oracle Performance*<sup>11</sup> for even more.)

Why is the knee value so important? For systems with randomly timed service requests, allowing

## Open Debate about Knees

In this article, I write about knees in performance curves, their relevance, and their application. Whether it's even worthwhile to try to define the concept of *knee*, however, has been the subject of debate going back at least 20 years.

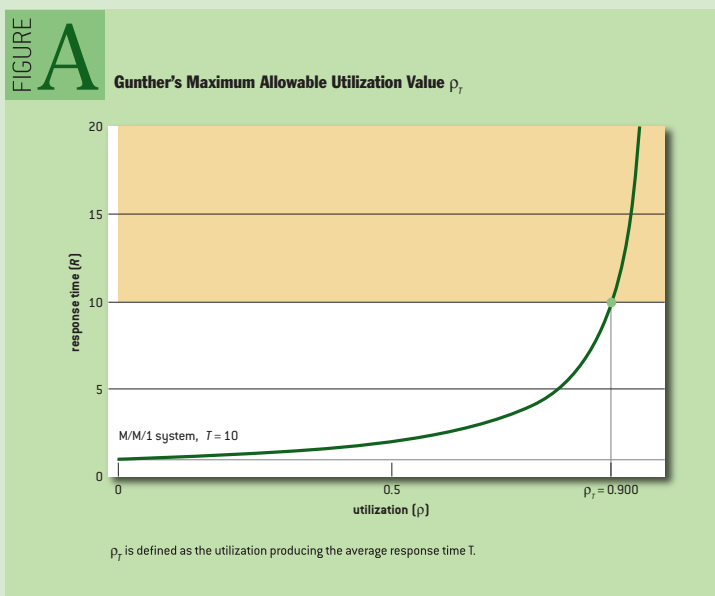
There is significant historical basis for the idea that the thing I've described as a knee in fact isn't really meaningful. In 1988, Stephen Samson argued that, at least for M/M/1 queuing systems, no "knee" appears in the performance curve. "The choice of a guideline number is not easy, but the rule-of-thumb makers go right on. In most cases there is not a knee, no matter how much we wish to find one," wrote Samson.<sup>3a</sup>

The whole problem reminds me, as I wrote in 1999,<sup>2a</sup> of the parable of the frog and the boiling water. The story says that if you drop a frog into a pan of boiling water, he will escape. But if you put a frog into a pan of cool water and slowly heat it, then the frog will sit patiently in place until he is boiled to death.

With utilization, just as with boiling water, there is clearly a "death zone," a range of values in which you can't afford to run a system with random arrivals. But where is the border of the death zone? If you are trying to implement a procedural approach to managing utilization, you need to know.

My friend Neil Gunther (see [http://en.wikipedia.org/wiki/Neil\\_J.\\_Gunther](http://en.wikipedia.org/wiki/Neil_J._Gunther) for more information about Neil) has debated with me privately that, first, the term *knee* is completely the wrong word to use here, in the absence of a functional discontinuity. Second, he asserts that the boundary value of .5 for an M/M/1 system is

wastefully low, that you ought to be able to run such a system successfully at a much higher utilization value than that. Finally, he argues that any such special utilization value should be defined expressly as the utilization value beyond which your average response time exceeds your tolerance for average response time (figure A). Thus, Gunther argues that any useful not-to-exceed utilization value is derivable only from inquiries about human preferences, not from



sustained resource loads in excess of the knee value results in response times and throughputs that will fluctuate severely with microscopic changes in load. Hence, on systems with random request arrivals, it is vital to manage load so that it will not exceed the knee value.

mathematics. (See [http://www.cmg.org/measureit/issues/mit62/m\\_62\\_15.html](http://www.cmg.org/measureit/issues/mit62/m_62_15.html) for more information about his argument.)

The problem I see with this argument is illustrated in figure B. Imagine that your tolerance for average response time is  $T$ , which creates a maximum tolerated utilization value of  $\rho_T$ . Notice that even a tiny fluctuation in average utilization near  $\rho_T$  will result in a huge fluctuation in average response time.

I believe that your customers feel the variance, not the mean. Perhaps they say they will accept average response times up to  $T$ , but humans will not be tolerant of performance on a system when a 1 percent change in average utilization over a one-minute period results in, say, a tenfold increase in average response time over that period.

I do understand the perspective that the knee values I've listed in this article are below the utilization values that many people feel safe in exceeding, especially for lower-order systems such as M/M/1. It is important, however, to avoid running resources at average utilization values where small fluctuations in utilization yield too-large fluctuations in response time.

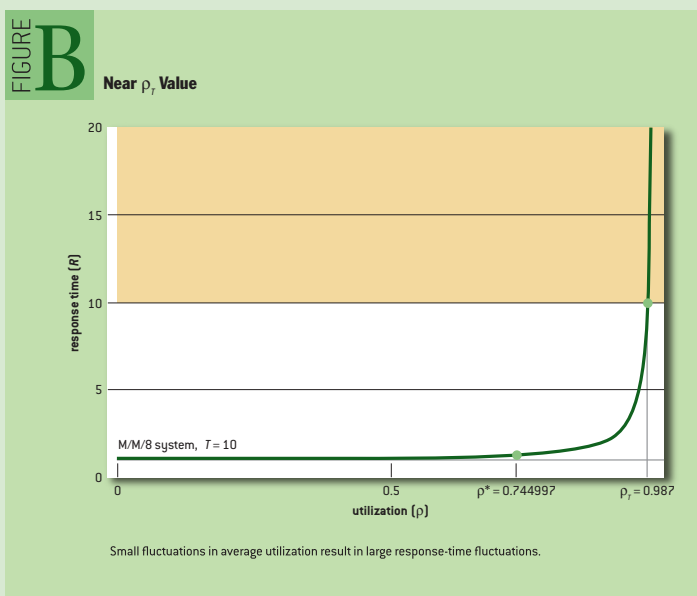
Having said that, I don't yet have a good definition for a *too-large fluctuation*. Perhaps, like response-time tolerances, different people have different tolerances for fluctuation. But perhaps there is a fluctuation tolerance factor that applies with reasonable universality across all users. The Apdex Application Performance Index standard, for example, assumes that the response time  $F$  at which users become "frustrated" is universally four times the response time  $T$  at which their attitude shifts from being "satisfied" to merely "tolerating."<sup>1a</sup>

The knee, regardless of how you define it or what we end up calling it, is an important parameter to the capacity-planning procedure that I described earlier in the main text of this article, and I believe it is an important parameter to the daily process of computer system workload management.

I will keep studying.

#### REFERENCES

- 1a. Apdex; <http://www.apdex.org>.
- 2a. Millsap, C. 1999. Performance management: myths and facts; <http://method-r.com>.
- 3a. Samson, S. 1988. MVS performance legends. In *Computer Measurement Group Conference Proceedings*: 148–159.



## RELEVANCE OF THE KNEE

How important can this knee concept be, really? After all, as I've told you, the M/M/m model assumes this ridiculously utopian idea that the system you're thinking about scales perfectly. I know what you're thinking: it doesn't.

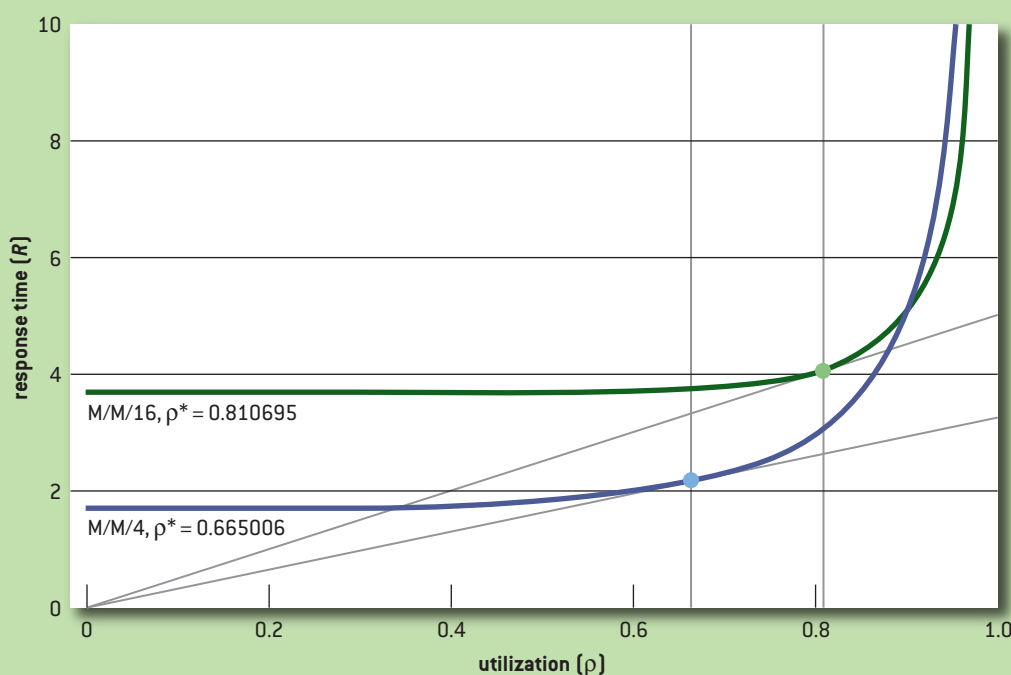
What M/M/m does give us is the knowledge that even if your system did scale perfectly, you would still be stricken with massive performance problems once your average load exceeded the knee values in table 1. Your system isn't as good as the theoretical systems that M/M/m models. Therefore, the utilization values at which *your* system's knees occur will be more constraining than the values in table 1. (I use the plural of *values* and *knees*, because you can model your CPUs with one model, your disks with another, your I/O controllers with another, and so on.)

To recap:

- Each of the resources in your system has a knee.
- That knee for each of your resources is less than or equal to the knee value you can look up in table 1. The more imperfectly your system scales, the smaller (worse) your knee value will be.
- On a system with random request arrivals, if you allow your sustained utilization for any resource on your system to exceed your knee value for that resource, then you'll have performance problems.
- Therefore, it is vital that you manage your load so that your resource utilizations will not exceed your knee values.

## FIGURE 5

Response Time Curves Showing Knees



The knee occurs at the utilization at which a line through the origin is tangent to the response time curve.



## CAPACITY PLANNING

Understanding the knee can collapse a lot of complexity out of your capacity planning. It works like this:

- Your goal capacity for a given resource is the amount at which you can comfortably run your tasks at peak times without driving utilizations beyond your knees.
- If you keep your utilizations less than your knees, your system behaves roughly linearly—no big hyperbolic surprises.
- If you're letting your system run any of its resources beyond their knee utilizations, however, then you have performance problems (whether you're aware of them or not).
- If you have performance problems, then you don't need to be spending your time with mathematical models; you need to be spending your time fixing those problems by rescheduling load, eliminating load, or increasing capacity.

That's how capacity planning fits into the performance management process.

## RANDOM ARRIVALS

You might have noticed that I used the term *random arrivals* several times. Why is that important?

Some systems have something that you probably don't have right now: completely deterministic job scheduling. Some systems—though rare these days—are configured to allow service requests to enter the system in absolute robotic fashion, say, at a pace of one task per second. And by this, I don't mean at an average rate of one task per second (for example, two tasks in one second and zero tasks in the next); I mean one task per second, as a robot might feed car parts into a bin on an assembly line.

If arrivals into your system behave completely deterministically—meaning that you know *exactly* when the next service request is coming—then you can run resource utilizations beyond their knee utilizations without necessarily creating a performance problem. On a system with deterministic arrivals, your goal is to run resource utilizations up to 100 percent without cramming so much workload into the system that requests begin to queue.

The reason the knee value is so important on a system with *random* arrivals is that these tend to cluster and cause temporary spikes in utilization. These spikes need enough spare capacity to

**TABLE 5**  
**M/M/m Knee Values for  
Common Values of m**

Service channel count	Knee utilization
1	50%
2	57%
4	66%
8	74%
16	81%
32	86%
64	89%
128	92%

consume so that users don't have to endure noticeable queuing delays (which cause noticeable fluctuations in response times) every time a spike occurs.

Temporary spikes in utilization beyond your knee value for a given resource are OK as long as they don't exceed a few seconds in duration. How many seconds are too many? I believe (but have not yet tried to prove) that you should at least ensure that your spike durations do not exceed eight seconds. (You'll recognize this number if you've heard of the "eight-second rule."<sup>2</sup>) The answer is certainly that if you're unable to meet your percentile-based response time promises or your throughput promises to your users, then your spikes are too long.

#### COHERENCY DELAY

Your system doesn't have theoretically perfect scalability. Even if I've never studied your system specifically, it's a pretty good bet that no matter what computer application system you're thinking of right now, it does *not* meet the M/M/m "theoretically perfect scalability" assumption. *Coherency delay* is the factor that you can use to model the imperfection.<sup>5</sup> It is the duration that a task spends communicating and coordinating access to a shared resource. Like response time, service time, and queuing delay, coherency delay is measured in time per task execution, as in seconds per click.

I won't describe a mathematical model for predicting coherency delay, but the good news is that if you profile your software task executions, you'll see it when it occurs. In Oracle, timed events such as the following are examples of coherency delay:

- enqueue
- buffer busy waits
- latch free

You can't model such coherency delays with M/M/m. That's because M/M/m assumes that all *m* of your service channels are parallel, homogeneous, and independent. That means the model assumes that after you wait politely in a FIFO queue for long enough that all the requests that enqueued ahead of you have exited the queue for service, it will be your turn to be serviced. Coherency delays don't work like that, however.

**EXAMPLE 10.** Imagine an HTML data-entry form in which one button labeled "Update" executes a SQL update statement, and another button labeled "Save" executes a SQL commit statement. An application built like this would almost guarantee abysmal performance. That's because the design makes it possible—quite likely, actually—for a user to click Update, look at his calendar, realize "uh-oh, I'm late for lunch," and then go to lunch for two hours before clicking Save later that afternoon.

The impact to other tasks on this system that wanted to update the same row would be devastating. Each task would necessarily wait for a lock on the row (or, on some systems, worse: a lock on the row's page) until the locking user decided to go ahead and click Save—or until a database administrator killed the user's session, which of course would have unsavory side effects to the person who thought he had updated a row.

In this case, the amount of time a task would wait for the lock to be released would have nothing to do with how busy the system was. It would be dependent upon random factors that exist outside of the system's various resource utilizations. That's why you can't model this kind of thing in M/M/m, and it's why you can never assume that a performance test executed in a unit-testing type of environment is sufficient for a making a go/no-go decision about insertion of new code into a production system.

## PERFORMANCE TESTING

All this talk of queuing delays and coherency delays leads to a very difficult question: How can you possibly test a new application enough to be confident that you're not going to wreck your production implementation with performance problems?

You can model. And you can test.<sup>1</sup> Nothing you do will be perfect, however. It is extremely difficult to create models and tests in which you'll foresee all your production problems in advance of actually encountering those problems in production.

Some people allow the apparent futility of this observation to justify not testing at all. Don't get trapped in that mentality. The following points are certain:

- You'll catch a lot more problems if you try to catch them prior to production than if you don't even try.
- You'll never catch all your problems in preproduction testing. That's why you need a reliable and efficient method for solving the problems that leak through your preproduction testing processes.

Somewhere in the middle between "no testing" and "complete production emulation" is the right amount of testing. The right amount of testing for aircraft manufacturers is probably more than the right amount of testing for companies that sell baseball caps. But don't skip performance testing altogether. At the very least, your performance-test plan will make you a more competent diagnostician (and clearer thinker) when the time comes to fix the performance problems that will inevitably occur during production operation.

## MEASURING

People feel throughput and response time. Throughput is usually easy to measure, response time is much more difficult. (Remember, throughput and response time are *not* reciprocals.) It may not be difficult to time an end-user action with a stopwatch, but it might be very difficult to get what you really need, which is the ability to drill down into the details of why a given response time is as large as it is.

Unfortunately, people tend to measure what's easy to measure, which is not necessarily what they *should* be measuring. It's a bug. Measures that aren't what you need, but that are easy enough to obtain and seem related to what you need are called *surrogate measures*. Examples include subroutine call counts and samples of subroutine call execution durations.

I'm ashamed that I don't have greater command of my native language than to say it this way, but here is a catchy, modern way to express what I think about surrogate measures: *surrogate measures suck*.

Here, unfortunately, *suck* doesn't mean *never work*. It would actually be better if surrogate measures never worked. Then nobody would use them. The problem is that surrogate measures work *sometimes*. This inspires people's confidence that the measures they're using should work all the time, and then they don't. Surrogate measures have two big problems.

- They can tell you your system's OK when it's not. That's what statisticians call *type I error*, the false positive.
- They can tell you that something is a problem when it's not. That's a *type II error*, the false negative.

I've seen each type of error waste years of people's time.

When the time comes to assess the specifics of a real system, your success is at the mercy of how good the measurements are that your system allows you to obtain. I've been fortunate to work in the

Oracle market segment, where the software vendor at the center of our universe participates actively in making it possible to measure systems the right way. Getting application software developers to use the tools that Oracle offers is another story, but at least the capabilities are there in the product.

#### PERFORMANCE IS A FEATURE

Performance is a software application feature, just like recognizing that it's convenient for a string of the form "Case 1234" to automatically hyperlink over to case 1234 in your bug-tracking system. (FogBugz, which is software that I enjoy using, does this.) Performance, like any other feature, doesn't just happen; it has to be designed and built. To do performance well, you have to think about it, study it, write extra code for it, test it, and support it.

Like many other features, however, you can't know exactly how performance is going to work out while you're still writing, studying, designing, and creating the application. For many applications (arguably, for the vast majority), performance is completely unknown until the production phase of the software development life cycle. What this leaves you with is this: since you can't know how your application is going to perform in production, you need to write your application so that it's easy to *fix* performance in production.

As David Garvin has taught us, it's much easier to manage something that's easy to measure.<sup>3</sup> Writing an application that's easy to fix in production begins with an application that's easy to measure in production.

Usually, when I mention the concept of production performance measurement, people drift into a state of worry about the measurement-intrusion effect of performance instrumentation. They immediately enter a mode of data-collection compromise, leaving only surrogate measures on the table. Won't software with an extra code path to measure timings be slower than the same software without that extra code path?

I like an answer that I once heard Tom Kyte give in response to this question.<sup>7</sup> He estimated that the measurement-intrusion effect of Oracle's extensive performance instrumentation is -10 percent or less (where *or less* means *or better*, as in -20 percent, -30 percent, etc.). He went on to explain to a now-vexed questioner that the product is at least 10 percent faster now because of the knowledge that Oracle Corporation has gained from its performance instrumentation code, more than making up for any "overhead" the instrumentation might have caused.

I think that vendors tend to spend too much time worrying about how to make their measurement code path efficient without figuring out first how to make it effective. It lands squarely upon the idea that Knuth wrote about in 1974 when he said that "premature optimization is the root of all evil."<sup>6</sup> The software designer who integrates performance measurement into a product is much more likely to create a fast application and—more importantly—one that will become faster over time.

#### REFERENCES

1. CMG (Computer Measurement Group, a network of professionals who study these problems very, very seriously); <http://www.cmg.org>.
2. Eight-second rule; [http://en.wikipedia.org/wiki/Network\\_performance#8-second\\_rule](http://en.wikipedia.org/wiki/Network_performance#8-second_rule).
3. Garvin, D. 1993. Building a learning organization. *Harvard Business Review* (July).

4. General Electric Company. What is Six Sigma? The roadmap to customer impact. <http://www.ge.com/sixsigma/SixSigma.pdf>.
5. Gunther, N. 1993. Universal Law of Computational Scalability; [http://en.wikipedia.org/wiki/Neil\\_J.\\_Gunther#Universal\\_Law\\_of\\_Computational\\_Scalability](http://en.wikipedia.org/wiki/Neil_J._Gunther#Universal_Law_of_Computational_Scalability).
6. Knuth, D. 1974. Structured programming with Go To statements. *ACM Computing Surveys* 6(4): 268.
7. Kyte, T. 2009. A couple of links and an advert...; <http://tkyte.blogspot.com/2009/02/couple-of-links-and-advert.html>.
8. Millsap, C. 2009. My whole system is slow. Now what? <http://carymillsap.blogspot.com/2009/12/my-whole-system-is-slow-now-what.html>.
9. Millsap, C. 2009. On the importance of diagnosing before resolving. <http://carymillsap.blogspot.com/2009/09/on-importance-of-diagnosing-before.html>.
10. Millsap, C. 2009. Performance optimization with Global Entry. Or not? <http://carymillsap.blogspot.com/2009/11/performance-optimization-with-global.html>.
11. Millsap, C., Holt, J. 2003. *Optimizing Oracle Performance*. Sebastopol, CA: O'Reilly.
12. Oak Table Network; <http://www.oaktable.net>.

#### ACKNOWLEDGMENTS

Thank you, Baron Schwartz, for the e-mail conversation in which you thought I was helping you, but in actual fact, you were helping me come to grips with the need for introducing coherency delay more prominently into my thinking. Thank you, Jeff Holt, Ron Crisco, Ken Ferlita, and Harold Palacio, for the daily work that keeps the company going and for the lunchtime conversations that keep my imagination going. Thank you, Tom Kyte, for your continued inspiration and support. Thank you, Mark Farnham, for your helpful suggestions. And thank you, Neil Gunther, for your patience and generosity in our ongoing discussions about knees.

**CARY MILLSAP** is the founder and president of Method R Corporation (<http://method-r.com>), a company devoted to software performance. He is the author (with Jeff Holt) of *Optimizing Oracle Performance* (O'Reilly) and a co-author of *Oracle Insights: Tales of the Oak Table* (Apress). He is the former vice president of Oracle Corporation's System Performance Group and a co-founder of Hotsos. He is also an Oracle ACE Director and a founding partner of the Oak Table Network, an informal association of well-known "Oracle scientists." He blogs at <http://carymillsap.blogspot.com>, and he tweets at <http://twitter.com/CaryMillsap>.